1 Speech Recognition and Signal Analysis by Exact Fast

2 Search of Subsequences with Maximal Confidence

Measure

3

4 SPECIFICATION

## 5 1 TITLE OF THE INVENTION

6 Speech Recognition and Signal Analysis by Exact Fast Search of Subsequences with Maximal

7 Confidence Measure

## 8 2 REFERENCE TO APPENDIX SUBMITTED ON CD

9 Not Applicable

## 10 3 CROSS-REFERENCE TO RELATED APPLICATION

11 This patent application has as parent application the patent application C99-00214/25.02.1999

12 registered with the State Office for Inventions and Trademarks (OSIM) in Bucharest, Ro-

13 mania. The present application is the US national stage of the international application

14 PCT/IB00/00189 registered with the International Patent Office in Geneva.

# 4    BACKGROUND OF THE INVENTION

15

## 4.1    FIELD OF THE INVENTION

16

17    The invention relates to a common component of:

18    • Speech Recognition, more particularly to the fields of Keyword Spotting and decoding,

19    • Segments Alignment for DNA and proteins,

20    • Recognition of Objects in Images,

## 4.2    DESCRIPTION OF THE RELATED ART

21

22    This invention addresses the problem of *keyword spotting (KWS)* in unconstrained speech

23    without explicit modeling of non-keyword segments (typically done by using filler HMM

24    models or an ergodic HMM composed of context dependent or independent phone models

25    without lexical constraints). Several methods (sometimes referred to as "sliding model meth-

26    ods") tackling this type of problem have already been proposed in the past. E.g., they use

27    Dynamic Time Warping (DTW) or Viterbi matching allowing relaxation of the (begin and

28    endpoint) constraints. These are known to require the use of an "appropriate" normaliza-

29    tion of the matching scores since segments of different lengths have then to be compared.

30    However, given this normalization and the relaxation of begin/endpoints, straightforward

31    Dynamic Programming (DP) is no longer optimal (or, in other words, the DP optimality

32    principle is no longer valid) and has to be adapted, involving more memory and CPU. In-

33    deed, at any possible ending time $e$, the match score of the best warp and start time $b$ of

34    the reference has to be computed (for all possible start times $b$ associated with unpruned

35 paths). Finally, this adapted DP quickly becomes even more complex (or intractable) for

36 more advanced scoring criteria (such as the confidence measures mentioned below).

37      Work in the field of confidence level, and in the framework of hybrid HMM/ANN systems

38 has shown that the use of accumulated local posterior probabilities (as obtained at the

39 output of a multilayer perceptron) normalized by the length of the word segment (or, better,

40 involving a double normalization over the number of phones and the number of acoustic

41 frames in each phone) was yielding good confidence measures and good scores for the re-

42 estimation of $N$-best hypotheses. However, so far the evaluation of such confidence measures

43 involved the estimation and rescoring of N-best hypotheses.

44      KWS methods without filler models have in common the selection of a subsequence of

45 the utterance to match the interesting keyword models. Let $X = \{x_1, x_2, \ldots, x_n, \ldots, x_N\}$

46 denote the sequence of acoustic vectors in which we want to detect a keyword, and let $M$

47 be the HMM model of a keyword $M$ and consisting of $L$ states $\mathcal{Q} = \{q_1, q_2, \ldots, q_\ell, \ldots, q_L\}$.

48 Assuming that $M$ is matched to a subsequence $X_b^e = \{x_b, \ldots, x_e\}$ $(1 \leq b \leq e \leq N)$ of $X$,

49 and that we have an implicit (not modeled) *garbage/filler state* $q_G$ preceding and following

50 $M$, one can define (approximate) the log posterior of a model $M$ given a subsequence $X_b^e$ as

51 the average posterior probability along the optimal path, i.e.:

$$-\log P(M|X_b^e) \simeq \frac{1}{e-b+1} \min_{\forall Q \in M} -\log P(Q|X_b^e)$$

$$\simeq \frac{1}{e-b+1} \min_{\forall Q \in M} \{-\log P(q^b|q_G)$$

$$-\sum_{n=b}^{e-1}[\log P(q^n|x_n) + \log P(q^{n+1}|q^n)]$$

$$-\log P(q^e|x_e) - \log P(q_G|q^e)\} \qquad (1)$$

56 where $Q = \{q^b, q^{b+1}, \ldots, q^e\}$ represents one of the possible paths of length $(e-b+1)$ in $M$, and

3

57 $q^n$ the HMM state visited at time $n$ along $Q$, with $q^n \in \mathcal{Q}$. In this expression, $q_G$ represents

58 the "garbage" (filler) state which is simply used here as the non-emitting initial and final

59 state of $M$. Transition probabilities $P(q^b|q_G)$ and $P(q_G|q^e)$ can be interpreted as the keyword

60 entrance and exit penalties, but can be simply set to 1. Local posteriors $P(q_\ell|x_n)$ can be

61 estimated using any of the known techniques: multi-gaussians, code-books, or as output

62 values of a multilayer perceptron (MLP) used in hybrid HMM/ANN systems. For a specific

63 sub-sequence $X_b^e$, expression (1) can easily be estimated by dynamic programming since the

64 sub-sequence and the associated normalizing factor $(e - b + 1)$ are given. However, in the

65 case of keyword spotting, this expression should be estimated for all possible begin/endpoint

66 pairs $\{b, e\}$ (as well as for all possible word models), and we define the matching score of $X$

67 on $M$ as:

$$S(M|X) = -\log P(M|X_{b^*}^{e^*}) \tag{2}$$

69 where the optimal begin/endpoints $\{b^*, e^*\}$, and the associated optimal path $Q^*$, are the

70 ones yielding the lowest average local posterior:

$$\langle Q^*, b^*, e^* \rangle = \operatorname*{argmin}_{\{Q,b,e\}} \frac{-1}{e - b + 1} \log P(Q|X_b^e) \tag{3}$$

72 Of course, in the case of several keywords, all possible models will have to be evaluated.

73 A double averaging involving the number of frames per phone and the number of phones

74 usually yields slightly better performance when used to rescore N-best candidates:

$$\langle Q^*, b^*, e^* \rangle = \tag{4}$$

$$\operatorname*{argmin}_{\{Q,b,e\}} \frac{-1}{J} \sum_{j=1}^{J} \left( \frac{1}{e_j - b_j + 1} \sum_{n=b_j}^{e_j} \log P(q_j^n|x_n) \right) nonumber \tag{5}$$

77 where $J$ represents the number of phones in the hypothesized keyword model and $q_j^n$ the

78  hypothesized phone $q_j$ for input frame $x_n$. However, given the time normalization and

79  the relaxation of begin/endpoints, straightforward DP is no longer optimal and has to be

80  adapted, usually involving more memory and CPU.

81  Filler-based KWS need a simpler decoding step. Although various solutions have been

82  proposed towards the direct optimization of (2), most of the keyword spotting approaches

83  today prefer to preserve the optimality and simplicity of Viterbi DP by modeling the complete

84  input and explicitly or implicitly modeling non-keyword segments by using so called filler or

85  garbage models as additional reference models. In this case, we assume that non-keyword

86  segments are modeled by extraneous garbage models/states $q_G$ (and grammatical constraints

87  ruling the possible keyword/non-keyword sequences).

88  [It is sufficient to consider only the case of detecting one keyword] _Let

89  *us consider only the case of detecting one keyword_* per utterance at a time. In this case,

90  the keyword spotting problem amounts at matching the whole sequence $X$ of length $N$ onto

91  an extended HMM model $\overline{M}$ consisting of the states $\{q_G, q_1, \ldots, q_L, q_G\}$, in which a path

92  (of length $N$) is denoted $\overline{Q} = \{\overbrace{q_G, \ldots q_G}^{b-1}, q^b, q^{b+1}, \ldots, q^e, \overbrace{q_G, \ldots q_G}^{N-e}\}$ with $(b-1)$ garbage states

93  $q_G$ preceding $q^b$ and $(N-e)$ states $q_G$ following $q^e$, and respectively emitting the vector

94  sequences $X_1^{b-1}$ and $X_{e+1}^N$ associated with the non-keyword segments.

95  Given some estimation of $P(q_G|x_n)$ (e.g., using probability density functions trained on

96  non keyword utterances), the optimal path $\overline{Q^*}$ (and, consequently $b^*$ and $e^*$) is then given

97  by:

98
$$\overline{Q^*} = \underset{\forall \overline{Q} \in \overline{M}}{\text{argmin}} - \log P(\overline{Q}|X)$$

99
$$= \underset{\forall \overline{Q} \in \overline{M}}{\text{argmin}} \{ - \log P(Q|X_b^e)$$

5

$$-\sum_{n=1}^{b-1} \log P(q_G|x_n) - \sum_{n=e+1}^{N} \log P(q_G|x_n)\}\qquad(6)$$

which can be solved by straightforward DP (since all paths have the same length). The main

problem of filler-based keyword spotting approaches is then to find ways to best estimate

$P(q_G|x_n)$ in order to minimize the error introduced by the approximations. Sometimes this

value was defined as the average of the $N$ best local scores while, in other approaches, this

value is generated from explicit filler HMMs. However, these approaches will usually not

lead to the "optimal" solution given by (2).

## 5  BRIEF SUMMARY OF THE INVENTION

The invention belongs to the technical domain of decoding, classification, alignment and

matching of data.

The invention introduces a new method performing tasks in keyword spotting in utter-

ances, detection of subsequences in chains of organic matter (DNA and proteins) and recog-

nition of objects in images. The proposed methods search in an optimized way the matching

that maximizes, over all the possible matchings, certain confidence measures based on nor-

malized posteriors. Three such confidence measures are used, two existed in previous work

in Speech Recognition, and the third one is a new one.

Application fields for this invention are: man-machine interfaces (using speech recogni-

tion; ex: control systems, banking, flight services, etc), coordination systems (for industrial

robots and automata) and development systems for pharmaceutic products.

## 6 BRIEF DESCRIPTION OF THE SEVERAL VIEWS OF THE DRAWINGS

Not Applicable

## 7 DETAILED DESCRIPTION OF THE INVENTION

[The present invention introduces a fast iterative method,] = *In the following, we show that it is possible to define an iterative process,* = referred to as *Iterating Viterbi Decoding (IVD)* with good/fast convergence properties, estimating the value of $P(q_G|x_n)$ such that straightforward DP (6) yields exactly the same segmentation (and recognition results) than (3). While the same result could be achieved through a modified DP in which all possible combinations (all possible begin/endpoints) would be taken into account, the method proposed below is much more efficient (in terms of both CPU and memory requirements).

Compared to previously devised "sliding model" methods the first method proposed here is based on:

1. A matching score defined as the average observation probability (posterior) along the most likely state sequence. It is indeed believed that local posteriors are more appropriate to the task.

2. The iteration of a Viterbi decoding algorithm, which does not require scoring for all begin/endpoints or N-best rescoring, and which can be proved to (quickly) converge to the "optimal" (from the point of view of the chosen scoring functions) solution without

138      requiring any specific filler models, using straightforward Viterbi alignments (similar

139      to regular filler-based KWS, but for some versions at the cost of a few iterations).

140     The IVD method is based on a similar criterion as the filler based approaches (6), but

141 rather than looking for explicit (and empirical) estimates of $P(q_G|x_n)$ we aim at mathe-

142 matically estimating its value (which will be different and adapted to each utterance) such

143 that solving (6) is equivalent to solving (3). Thus, we perform an iterative estimation of

144 $P(q_G|x_n)$, such that the segmentation resulting of (6) is the same than what would be ob-

145 tained from (3). Defining $\varepsilon_t = -\log P(q_G|x_n)$ at iteration $t$, the proposed method can be

146 summarized as follows:

147     1. Start the first iteration, $t = 0$, from an initial value $\varepsilon_0 = \Pi$ (it is actually proven that

148        the iterative process presented here will always converge to the same solution, in more

149        or less cycles with the worst case upper bound of N iterations, independently of this

150        initialization, e.g., with $\Pi$ equal with a cheap estimation of the score of a "match").

151        In one of the developed versions, $\varepsilon_0$ is initialized to $-\log$ of the maximum of the local

152        probabilities $P(q_k|x_n)$ for each frame $x_n$.

153        An alternative choice is to initialize $\varepsilon_0$ to a pre-defined threshold score, T, that expres-

154        sion (1) should reach to declare a keyword "matching" (see step 4 below). In this last

155        case, if $\varepsilon_1 > \varepsilon_0$ at the first iteration, then we can (as proven) directly infer that the

156        match will be rejected, otherwise it will be accepted.

157     2. Given the estimate $\varepsilon_t$ of $P(q_G|x_n)$ at current iteration $t$, find the optimal path $\langle \overline{Q}_t, b_t, e_t \rangle$

158        according to (6) and matching the complete input.

159    3. Estimate the value of $\varepsilon_{t+1}$ to be used in the next iteration as the average of the local

160       posteriors along the optimal path $Q_t$ (matching the $X_{b_t}^{e_t}$ resulting of (6) on the keyword

161       model) i.e.:

$$\varepsilon_{t+1} = -\frac{1}{(e_t - b_t + 1)} \log P(Q_t | X_{b_t}^{e_t}) \tag{7}$$

162

163    4. Increment $t$ and return to (2) iterating until convergence is detected. If we are not

164       interested in the optimal segmentation, this process could also be stopped as soon as it

165       reaches a $\varepsilon_{t+1}$ lower than a (pre-defined) minimum threshold, T, below which we can

166       declare that a keyword has been detected.

167    Correctness and convergence proof of this process and generalization to other criteria, are

168    available: each IVD iteration (from the second iteration) will decrease the value of $\varepsilon_t$, and the

169    final path yields the same solution than (3). The above method has a very good experimental

170    convergence speed (3-5 iterations in our tests). For one version of IVD (when $\varepsilon_0$ is initialized

171    using the acceptance threshold, T), the detection is decided after one single step.

172       A version with the same effort but suboptimal results is proposed in the following para-

173    graph. Let $T(\overline{M}, X)$ be a matrix holding the HMM emission probabilities for an utterance

174    $X$ whose time-frames define the columns, and where the states of the hypothesized word

175    $W$ define the rows. When using the standard DP, one computes for each element of the

176    matrix $T(\overline{M}, X)$ at frame $k$ of $X$ and state $s$ of $\overline{M}$ three values: $S_{ks}$, $L_{ks}$ and $C_{ks}$, where

177    $S_{ks}$ corresponds to the sum of the entries on the optimal path that leads to the entry, $L_{ks}$

178    holds the length of the optimal path computed so far, and $C_{ks}$ is the estimation of the cost

179    on the optimal expanded path. By a path leading to an entry $T(k, s)$ we mean a sequence

180    of entries in the table $T$, such that there is exactly an entry for each time frame $t \leq k$. At

181   each entry $T(k, s)$, DP selects a locally optimal path noted $P_{ks}$. At each step $k$, we consider

182   all pairs of entries of table $T(\overline{M}, X)$ of type $T(k, s)$, $T(k-1, t)$. We update for each such

183   pair, the current cost $C_{ks}$ (initially $\infty$), by comparing it with the alternative given by:

184

$$S_{ks} = S_{(k-1)t} - \log p(\dot{s}|x_k)p(s|t)$$

185

$$L_{ks} = L_{(k-1)t} + 1, \forall t > 0, t \leq L$$

186

$$C_{ks} = \frac{S_k}{L_k} \qquad\qquad (8)$$

187   wanting to have at step $k$ the path $P_{ks}$ from the paths $P_{(k-1)t}$ that minimizes $C_{NL}$. With

188   DP, one will choose the $P_{ks}$ with minimal $C_{ks}$.

189     This version can yield suboptimal results since the optimality principle is not respected

190   by the expression 8. The optimality principle of Dynamic Programming requires that the

191   path to the frame $k-1$ that minimizes $C_{NL}$, also minimizes $C_{ks}$ for an entry at frame $k$ of

192   table $T(\overline{M}, X)$.

193     Another technique that is suboptimal in time and/or quality is obtained from the previous

194   one adopting a beam-search approach and a set of safe prunings. The Dynamic Programming

195   can be viewed as a set of safe prunings that are applied at each entry of the DP table and

196   has the property that only one alternative is maintained. Dynamic Programming cannot be

197   used, since the principle of optimality is not respected. The following types of safe pruning

198   that can be done are introduced by the present invention. Within the current invention we

199   found a set of safe prunings as follows: we have proved that if at a frame $a$ we have two paths

200   $P_a'$ and $P_a''$ with $S_a'' < S_a'$ and $L_a' < L_a''$, then at no frame $c \geq a$ will a path $P_c''$ be forsaken for

201   a path $P_c'$ if $P_a' \subset P_c'$, $P_a'' \subset P_c''$ and $P_c' \backslash P_a' \equiv P_c'' \backslash P_a''$. We will note the order relation as $P_a'' \prec P_a'$.

202 We have further shown that a path P' may be safely discarded only when we know a lower

203 cost one, P".

$$P' \prec P'' \Rightarrow C_k' < C_k'' \qquad (9)$$

205 Thus, the method described in following method computes $S(M, X)$ and $Q^*$ from equa-

206 tion (3). By ordering the set of paths, according to Equation 9, we only need to check the

207 step (1.1) of the following method up to the eventual insertion place. The last paths are

208 candidates for pruning in step (1.2). In order for the pruning to be acceptable, we will prune

209 only paths that were too long on the last state. An additional counter for each path is

210 needed for storing the state length. This counter is reset when an entry from another row

211 is added and is incremented at each advance with a frame. The following steps detail this

212 method for a model W and an utterance X:

213     a) Initialize all elements of a matrix, SetOfPaths(1..N, 1..K), to $\emptyset$

214     b) For all frames from 1 to N, for all states from 1 to K, for all candidates $p_i$ in

215        SetOfPaths(frame-1, 1..K):

216          − For all $p_j$ in SetOfPaths[frame, state], if $p_i \prec p_j$ then delete $p_j$ (1.1), and if $p_j \prec p_i$

217            then continue step b) (1.2)

218          − Insert $p_i$ in SetOfPaths[frame, state]

219     c) Select SetOfPaths[frame, K] as the best of the candidates

220     The next method builds on the previous technique and is a fast procedure for maximizing

221 a more complex confidence measure that yields better results in practice. The corresponding

222 confidence measure is defined as:

$$\frac{1}{NVP} \sum_{h_i \in VP} \frac{\sum_{pst \in h_i} -\log(pst)}{length(h_i)} \tag{10}$$

224 where NVP stands for the *number of visited phonemes* and VP stands for the *set of visited*

225 *phonemes*. An average is computed over all posteriors *pst* of the emission probabilities for the

226 time frames matched to the visited phoneme $h_i$. The function $length(h_i)$ gives the number of

227 time frames matched against $h_i$. This method uses a breath first Beam Search algorithm. It

228 exploits a set of reduction rules and certain normalizations. For the state $q_G$, in this method,

229 the logarithm of the emission posterior is equal with zero. For each frame $e$ and for each

230 state $s$, the set of paths/probabilities of having the frame $e$ in the state $s$ is computed as

231 the first $\mathcal{N}$ maxima ($\mathcal{N}$ can be finite) of the confidence measure for all paths in HMM $\overline{M}$ of

232 length $e$ and ending in the state $s$. The paths that according to the reduction rules will loose

233 the final race when compared with another already known path, will be deleted as well. Let

234 us note $a_1$, $p_1$, $l_1$, respectively $a_2$, $p_2$ and $l_2$ the confidence measure for the previously visited

235 phonemes, the posterior in the current phoneme and the length in the current phoneme for

236 the path $Q_1$, respectively the path $Q_2$. The rules that can be used for the reduction of the

237 search space by discarding a path $Q_1$ for a path $Q_2$ are in this case any of the next ones:

238      1. $l_2 \geq l_1$, $A > 0$, $B \leq 0$ and $L_c^2 A + L_c B + C \geq 0$

239      2. $l_2 \geq l_1$, $A \geq 0$, $B \geq 0$ and $C \geq 0$

240      3. $l_2 \geq l_1$, $A \leq 0$, $C \geq 0$ and $L^2 A + LB + C \geq 0$

241      4. $l_2 \geq l_1$, $A = 0$, $B < 0$ and $LB + C \geq 0$

242 where $A = a_1 - a_2$, $B = (a_1 - a_2)(l_1 + l_2) + p_1 - p_2$, $C=(a_1 - a_2)l_1l_2 + p_1l_2 - p_2l_1$, $L =$

243 $L_{max} - \max\{l_1, l_2\}$, $L_c = -B/2A \geq 0$ and $L_{max}$ is the maximum acceptable length for a

244 phoneme. By discarding paths only if one of the above rules is satisfied, the optimum defined

245 by the confidence measure with double normalization can be guaranteed, if no phone may be

246 avoided by the HMM $M$. Any HMM may be decomposed in HMMs with this quality. The

247 4-th rule is included in the 3-rd and its test is useless if the last one was already checked.

248 The first test, $l_2 \geq l_1$ tells us if $Q_2$ has chances to eliminate $Q_1$, otherwise we will check

249 if $Q_1$ eliminates $Q_2$. These tests were inferred from the conditions of maintaining the final

250 maximal confidence measure while reduction takes place. In order to use the method of

251 double normalization without decomposing HMMs that skip some phonemes, the previous

252 rules are modified taking into account the number of visited phonemes for any path $F_1$

253 respectively $F_2$ and the number of phonemes that may follow the current state. A simplified

254 test can be:

255     • $l_2 \geq l_1$, $A \geq 0$, $p_1 \geq p_2$ respectively $F_2 \geq F_1$ for the HMMs that skips phonemes.

256 This test is weaker than the $2^{nd}$ reduction rule. For example a path is eliminated by a second

257 path if the first one has an inferior confidence measure (higher in value) for the the previous

258 phonemes, a shorter length and the minus of the logarithm of the cumulated posterior in

259 the current phoneme also inferior (higher in value) to that of the second one. An additional

260 confidence measure based on the maximal length, $L_{max}$, and on the maximum of the minus

261 of the logarithm of the cumulated and normalized posterior in phoneme, $P_{max}$, can be used

262 in order to limit the number of stored paths.

263     • $p > L_{max}P_{max}$ in any state

264     • $\frac{p}{l} > P_{max}$ at the output from a phoneme

265 where p and l are the values in the current phoneme for the minus of the logarithm of

266 cumulated posterior and for the length of the path that is discarded. These tests allow for

267 the elimination of the paths that are too long without being outstanding, respectively of

268 the paths with phonemes having unacceptable scores, otherwise compensated by very good

269 scores in other phonemes. If $\mathcal{N}$ is chosen equal with one, the aforementioned rules are no

270 longer needed, but always we propagate the path with the maximal current estimation of

271 the confidence measure. The obtained results are very good, even if the defined optimum is

272 guaranteed for this method only when $\mathcal{N}$ is bigger than the length of the sequence allowed

273 by $L_{max}$ or of the tested sequence. The same approach is valid for the simple normalization,

274 where the HMM for the searched word will be grouped into a single phoneme.

275     The present invention can exploit a newly designed a confidence measure, version named

276 "Real Fitting", that represents differently the exigencies of the recognition. Since the

277 phonemes and the absent states can be modeled by the used HMMs, we find it interest-

278 ing to request the fitting of each phoneme in the model with a section of the sequence.

279 Therefore, we measure the confidence level of a subsequence as being equal with the max-

280 imum over all phonemes of the minus of the logarithm of the cumulated posterior of the

281 phone, normalized with its length:

$$\max_{phonem \in Visited \ Phonems} \frac{\sum_{phonem} - \log(posteriors)}{phonem \ length} \tag{11}$$

282

283 The rule that may be used in this framework for the reduction of the number of visited paths

284 is:

285     • $Q_2$ is discarded in favor of another path $Q_1$ if the confidence measure of the Real

286         Fitting for the previous phonemes is inferior (higher in value) for $Q_2$ compared with

287         $Q_1$, and if $p_1 \leq p_2$ and $l_2 \leq l_1$.

288   where $p_1$, $l_1$, respectively $p_2$, $l_2$ represent the minus of the logarithm of the cumulated poste-

289   rior respectively the number of frames in the current phoneme for the path $Q_1$ respectively

290   $Q_2$. Similarly to the previous method, the set of visited paths can be pruned by discarding

291   those where:

292       • $p > L_{max} P_{max}$ in any state

293       • $\frac{p}{l} > P_{max}$ at the output from a phoneme

294   where p and l are the values in the current phoneme for the minus of the logarithm of the

295   cumulated posterior and for the length of the path that is discarded. We recall that the

296   meaning of the constants are the maximal length $L_{max}$, respectively the accepted maxima

297   of the minus of the logarithm of the cumulated and normalized posterior in phoneme, $P_{max}$.

298         This invention thus proposes a new method for keyword spotting, based on recent ad-

299   vances in confidence measures, using local posterior probabilities, but without requiring the

300   explicit use of filler models. A new method, referred to as *Iterating Viterbi Decoding (IVD)*,

301   to solve the above optimization problem with a simple DP process (not requiring to store

302   pointers and scores for all possible ending and start times). Other three new beam-search

303   algorithms corresponding to three different confidence measures are also proposed.

304         To summarize, the object of the invention consists of:

305       • Method of recognition of a subsequence using a direct maximization of confidence

306          measures.

307    • The method of IVD for directly maximizing the confidence measures based on simple

308      normalization.

309    • The use of the confidence measure and method of recognition named 'Real Fitting',

310      based on individual fitting for each phoneme.

311    • Methods of recognition using simple and double normalization by:

312    • combining these measures with additional confidence measures mentioned here, respec-

313      tively the maximal length and real matching limitation.

314    • The use of the aforementioned methods in keyword recognition.

315    • The use of the aforementioned methods in subsequence recognition of organic matter.

316    • The use of the aforementioned methods in recognition of objects in images.

317    DETAILED DESCRIPTION OF PREFERRED EMBODIMENTS

318    Execution: The method can be performed using a personal computer or can be imple-

319    mented in specialized hardware.

320    1. A representation under the form of an HMM is obtained for the subsequences that are

321       looked for (word, protein profile, section of an image of the object).

322    2. A tool will be obtained (eventually trained Ex: for speech recognition) for the esti-

323       mation of the posteriors. For example multi-Gaussians, neuronal networks, clusters,

324       database with Generalized Profiles and mutation matrices (PAM, BLOSSUM, etc.).

16

325  3. One of the proposed algorithms should be implemented. They yield close performance

326  but the method of Real Fitting coupled with a well checked dictionary should perform

327  best.

328  For the first algorithm (IVD)

329  (a) The classic algorithm of Viterbi is implemented with the modification that, for

330  each pair $P = \langle sample, state \rangle$ one propagates the time-frame of transition be-

331  tween the state $q_G$ and the states of the HMM $M$ for the path that arrives at P.

332  These are inherited from the path that wins the entrance in the pair P, excepting

333  for the moment when their decision is taken, namely when they receive the index

334  of the corresponding sample.

335  (b) $w = -\log P(M|X_t^e)$ is computed by subtracting from the cumulated posterior

336  that is returned by the Viterbi algorithm for the path $Q_{b_t}^{e_t}$, the value $(N - (e_t -$

337  $b_t + 1)) * \varepsilon_t$ corresponding to the contribution of the states $q_G$ and dividing the

338  result through $e_t - b_t + 1$. $e_t - b_t + 1$ from the previous formula can be factored

339  outside the fraction.

340  (c) The initialization of $\varepsilon$ is made with an expected mean value. One can use the $w$

341  that is computed when the state $q_G$ is associated with an emission posterior equal

342  to the average of the best $K$ emission probabilities of the current sample as done

343  in the well-known "garbage on-line model". In this case, K is trained using the

344  corresponding technique.

345  The next 'Beam search' algorithms, are implemented according to the description in

346   the corresponding sections. For each pair $P = \langle sample, state\rangle$ one computes for each

347   corresponding path the sum and length in the last phoneme, as well as the sum over

348   the normalized cumulated posteriors of the previous phonemes (and their number).

349   Also, the entrance and exit samples into the HMM $M$ are computed and propagated

350   like in the previous method, in order to ensure the localization of the subsequence.

351   4. If one searched entity (keyword, sequence, object) can have several HMM models, all

352   of them are taken into consideration as competitors. This is the case of the words

353   with several pronunciations (or of the objects that have different structures in different

354   states, for the recognition in images).

355   After the computation of the confidence measure for each model of the subsequences,

356   one eliminates those with a confidence measure in disagreement with a 'threshold' that

357   is trained for the configuration and the goal of the given application. For example, for

358   speech recognition with neuronal networks and minus of the logarithm of the posteriors,

359   the 'threshold' is chosen in the wanted point of the ROC curve obtained in tests.

360   5. The remained alternatives are extracted in the order of their confidence measure and

361   with the elimination of the conflicting alternatives until exhaustion. Each time when

362   an alternative is eliminated, the searched entity with the corresponding HMM is re-

363   estimated for the remaining sections in the sequence in which the search is performed.

364   If the new confidence measure passes the test of the 'threshold', then it will be inserted

365   in the position corresponding to its score in the queue of alternatives.

366   6. The successful alternatives can undergo tests of superior levels like for example a

367       question of confirmation for speech recognition, opinion of one operator, etc.

368   7. For objects recognition in images:

369       Posteriors are obtained by computing a distance between the color of the model and

370       that of element in the section of the image. If the context requires, the image will be

371       preprocessed to ensure a certain normalization (Ex: changeable conditions of light will

372       make necessary a transformation based on the histogram).

373       The phonemes of the speech recognition correspond to parts of the object. The struc-

374       ture (existence of transitions and their probabilities) can be modified, function of the

375       characteristics detected along the current path. For example, after detecting regions

376       of the object with certain lengths, one can estimate the expected length of the remain-

377       ing regions. Thus, the number of the expected samples for the future states can be

378       established and the HMM attached to the object will be configured accordingly.

379       A direction is scanned for the detection of the best fitting and afterwards, other direc-

380       tions will be scanned for discovering new fittings, as well as for testing the previous

381       ones. The final test will be certified by classical methods such as cross-correlation or

382       by the analysis of the contours in the hypothesized position.

383   To mention some examples for the application of the proposed method:

384       • The recognition of keywords begins to be used in answering automates of banking

385       system as well as telephone and automates for control, sales or information. The

386       method offers a possibility to recognize keywords in spontaneous speech with multiple

387       speakers.

388     • The recognition of DNA sequences is important for the study of the human Genome.

389       One of the biggest problem of the involved techniques consists in the high quantity of

390       data that have to be processed.

391     • The recognition of objects in images is used, among others, in cartography and in the

392       coordination of industrial robots. The method allows a quick estimation of the position

393       of the objects in scenes and can be validated with extra tests, using classical methods

394       of cross-correlation.

395 WE CLAIM:

396     1. (canceled) rewritten/re-presented in claim 5

397     2. (canceled) rewritten/re-presented in claim 6

398     3. (canceled) rewritten/re-presented in claim 7

399     4. (canceled) rewritten/re-presented in claim 8

400     5. (canceled) rewritten/re-presented in claim 9

401     6. (canceled) rewritten/re-presented in claim 10

402     7. (canceled) rewritten/re-presented in claim 11

403     8. (canceled) rewritten/re-presented in claim 12

404     9. (re-presented - formerly independent claim 5) A method of recognizing an observed

405     subsequence as being generated by one of a set of Hidden Markov Models (HMM),

406     characterized by:

407     • the fact that it searches the subsequence, $Q$, that offer the minimization of an

408     inverse confidence measure, over all possible matchings,

409     • where the inverse confidence measure is one of

410     1) the accumulated posterior, normalized with the length of the matched sub-

411     sequence $X_b^e$ (aka. 'simple normalization')

412
$$\frac{-1}{e-b+1} \log P(Q|X_b^e)$$

21

413      2) partitioning the states in a HMM into phonems, having a function

414         Phonemes(Q) that returns the segmentation of a path Q in the HMM into

415         phonems, and computing one of:

416         2a) the worst average match in a phoneme, called 'real fitting',

417
$$\operatorname*{argmin}_{Q} \left( \max_{Q \in Phonemes(Q)} \frac{\sum_{q^k \in Q} -\log P(q^k|x_k)}{|\{k|q^k \in Q\}|} \right)$$

418         2b) double normalization of the accumulated posterior over the number of

419         phonemes, J, and over the number of acoustic samples, $e_j - b_j + 1$, where

420         $e_j$ is the time frame where Q enters phoneme $j$, and $b_j$ is the exit time

421         frame from each phoneme, $j$,

422
$$\frac{-1}{J} \sum_{j=1}^{J} \left( \frac{1}{e_j - b_j + 1} \sum_{n=b_j}^{e_j} \log P(q_j^n|x_n) \right)$$

423    • and allows for the optional revaluation of the alternatives that offer the high-

424        est scores of a mentioned confidence measure on the basis of another confidence

425        measure,

426    • and when based on the confidence measure called 'simple normalization' uses a

427        method that applies Viterbi decoding for a HMM obtained by extending the initial

428        one with a filler state just after start and one just before the termination state,

429        and estimates the emission probability of the filler states in an iterative manner

430        as being equal to the inverse confidence measure in the previous iteration,

431        and where the emission probability in the filler states in the first iteration can be

432        initialized to any floating point number, but the iteration stops:

433       i at convergence yielding the estimation of a keyword's boundaries and score

434          as the obtained boundaries and score of non-filler states of the HMM,

435       ii when the confidence measure descends under a threshold value, T, estimating

436          only the keyword existence,

437       iii when the emission probability of filler states, $\varepsilon_0$ is initialized with T and is

438          reestimated, as value of $\varepsilon_1$ at the end of the first iteration, to be higher than

439          T deciding keyword inexistence,

440 • or for any of the three confidence measures: 'simple normalization', 'double nor-

441   malization' or 'real fitting', uses a beam-search-like algorithm that considers the

442   emission probability of the filler state as zero, computes progressively for each

443   pair of sample and state of HMM a set of possible alternatives paths to reach it,

444   the computation of this set is based on the sets of paths that lead to the states that

445   can be associated to the previous sample and extended with transitions allowed

446   by the analyzed HMM,

447   where this set can be reduced by using appropriate (safe) rules for the given

448   confidence measure, ensuring the correctness of the inference,

449   and where this set can be also reduced by using heuristics, for speeding up the

450   computation despite the risk of reducing the theoretical quality of the recognition,

451   heuristics of which a fast version stores only the best match,

452 and for all confidence measures one can prune the set of alternatives with safe rules

453 guaranteeing optimality, where:

454   • the 'simple normalization' confidence measure with beam-search is used with a

455     safe pruning that discards a path $Q_1$ given the existence of a path $Q_2$ whenever

456     $S_2 < S_1$ and $L_1 < L2$, where $S_1$ and $L_1$ respectively $S_2$ and $L_2$ are the minus of

457     the cumulated log of posteriors along the paths, and the lengths of the paths, for

458     the paths $Q_1$ respectively $Q_2$, and which can be optimized by sorting competing

459     paths based on their cost

460   • the 'double normalization' confidence measure on HMMs where no path skips any

461     phoneme is used with a safe pruning that discards a path $Q_1$ given the existence

462     of a path $Q_2$ whenever one of the following tests succeed:

463     (a) $l_2 \geq l_1$, $A > 0$, $B \leq 0$ and $L_c^2 A + L_c B + C \geq 0$

464     (b) $l_2 \geq l_1$, $A \geq 0$, $B \geq 0$ and $C \geq 0$

465     (c) $l_2 \geq l_1$, $A \leq 0$, $C \geq 0$ and $L^2 A + LB + C \geq 0$

466     (d) $l_2 \geq l_1$, $A = 0$, $B < 0$ and $LB + C \geq 0$

467     where we denote by $a_1$, $p_1$, $l_1$, respectively by $a_2$, $p_2$ and $l_2$ the confidence measure

468     for the previously visited phonemes, the posterior in the current phoneme and

469     the length in the current phoneme for the path $Q_1$, respectively the path $Q_2$,

470     and we also use the notations $A = a_1 - a_2$, $B = (a_1 - a_2)(l_1 + l_2) + p_1 - p_2$,

471     $C = (a_1 - a_2)l_1 l_2 + p_1 l_2 - p_2 l_1$, $L = L_{max} - \max\{l_1, l_2\}$, $L_c = -B/2A$ and $L_{max}$ is

472     the maximum acceptable length for a phoneme,

473   • the 'double normalization' confidence measure on HMMs where some paths skip

474     phonemes is used with a safe pruning that discards a path $Q_1$ given the existence

475     of a path $Q_2$ whenever $l_2 \geq l_1$, $A \geq 0$, $p_1 \geq p_2$ respectively $F_2 \geq F_1$,

476           where $F_1$ respectively $F_2$ are the number of visited phonemes for paths $Q_1$ and

477           $Q_2$,

478           • the 'real fitting' is used with the safe pruning: $Q_2$ is discarded in favor of another

479           path $Q_1$ if the confidence measure of the Real Fitting for the previous phonemes

480           is inferior (higher in value) for $Q_2$ compared with $Q_1$, and if $p_1 \leq p_2$ and $l_2 \leq l_1$,

481           where $p_1$, $l_1$, respectively $p_2$, $l_2$ represent the minus of the logarithm of the cumu-

482           lated posterior respectively the number of frames in the current phoneme for the

483           path $Q_1$ respectively $Q_2$,

484           • and besides the previously mentioned safe pruning, heuristic prunings are also

485           used for removing paths when $p > L_{max}P_{max}$ in any state or when $\frac{p}{l} > P_{max}$ at

486           the output from a phoneme,

487           where p and l are the values in the current phoneme for the minus of the logarithm

488           of cumulated posterior and for the length of the path that is discarded.

489   10. (re-presented - formerly dependent claim 6) The method of claim 9, where the method

490      is used to estimate the existence of keywords and their position in utterances, using

491      Hidden Markov Models that model keywords.

492   11. (re-presented - formerly dependent claim 7) The method of claim 9, where the method

493      is used to estimate the existence of biomolecular subsequences and their position in the

494      chains of DNA using hidden Markov models to model the searched subsequences, and

495      where these models can be obtained by trivial translation from generalized profiles.

496   12. (re-presented - formerly dependent claim 8) The method of claim 9, where it carries out

25

497     the estimation of the existence of objects and their position in images, characterized

498     by the fact that

499     • it uses models of objects as subsequences represented by Hidden Markov Models,

500     • namely sections through views of objects are modeled by Hidden Markov Models,

501     • it uses emission probabilities based on a distance computed between colors, sim-

502     ple distances being yield by a Gaussian with median at the target color, or a

503     normalized inverse of the Euclidean distance in the RGB space,

504     • wherein the Hidden Markov Models that model the objects can be structured of

505     distinct regions, that play in the frame of the method the role of the phonemes

506     in claim 9,

507     • and wherein the models of the objects can be modified in a dynamic manner during

508     decoding with respect to the transition properties (existence and probability) on

509     the basis of the so far accumulated information in the process.

# 510 8 ABSTRACT OF THE DISCLOSURE

511 The invention belongs to the technical domain of decoding, classification, alignment and 512 matching of data.

513 The invention introduces a new method performing tasks in keyword spotting in ut- 514 terances, detection of subsequences in chains of organic matter (DNA and proteins) and 515 recognition of objects in images. The proposed method searches in an optimized way the 516 matching that maximizes, over all the possible matchings, certain confidence measures based 517 on normalized posteriors. Three such confidence measures are used, two existed in previous 518 work in Speech Recognition, and the third one is a new one.

519 Application fields for this invention are: man-machine interfaces (using speech recogni- 520 tion; ex: control systems, banking, flight services, etc), coordination systems (for industrial 521 robots and automata) and development systems for pharmaceutic products.